

PATENT APPLICATION

**SYSTEM FOR DETERMINING DEGREES OF SIMILARITY IN
EMAIL MESSAGE INFORMATION**

Inventor: Tonny Yu, a citizen of the USA residing at
1710 Crockett Lane
Hillsborough, CA 94010

Assignee: MAILSHELL, INC.
695 Fifth St., #3
San Francisco, CA 94107

Entity: Small

SYSTEM FOR DETERMINING DEGREES OF SIMILARITY IN EMAIL MESSAGE INFORMATION

5

COPYRIGHT NOTICE

[01] A portion of the disclosure recited in the specification contains material which is subject to copyright protection. Specifically, a Source Code Appendix is provided that includes source code instructions for a process by which the present invention is practiced in a computer system. The copyright owner has no objection to the facsimile reproduction of the specification as filed in the Patent and Trademark Office. Otherwise all copyright rights are reserved.

BACKGROUND OF THE INVENTION

[02] Although electronic mail, or email, has become immensely popular and is a huge benefit for many users, today's email systems are also plagued by increasing numbers of unwanted mail, referred to as "spam." Spam email has reached such large proportions with respect to desired email that systems are now sought to defeat the sending and delivery of spam.

[03] For example, one approach is to design "filters" to block spam before it is received in a user's email in-box. The filters use different user-designed criteria such as detecting a sender's name, or a word or phrase used in a subject header. Filters can also be used to sort email into separate folders once the email has been received by a user so that the user can ignore folders into which spam is sorted. These approaches are not without shortcomings since the filters typically work on keyword matching or common and relatively easy to detect syntax or language features. Spam emailers have developed ways to thwart simple filter approaches. Sophisticated spam senders can use processes to modify an original email messages into different variations that each communicate essentially the same message. Typically the message is designed to sell something to a recipient, or is designed to provide other commercial advantage to the spam emailer.

[04] For example, one line in an email message might be "buy this now." The line can be modified to "you should try this now." Other properties of the message can be modified such as the order of sentences, addition or removal of words or phrases, changes in spacing or other message formatting, etc. Since the modified spam email messages are different, it is difficult for simple spam detection routines to successfully identify a primary characteristic of spam email, namely, that the email is sent in large number such as thousands, hundreds of thousands or more instances of the same message. Such high-volume email is referred to as "bulk" email. Spam emailers can also use such approaches to change other characteristics of an email message, such as sender identification, routing information and other information that may be associated with an email message that could otherwise help determine that the email message is a bulk emailing and is likely to be spam.

[05] Spam detection is further complicated because all bulk emailings are not necessarily spam. For example, if thousands of users desire to be informed of daily weather from a weather source then the messages are likely to be the same or similar, depending on the regional location of the users. Even though such email would qualify as bulk email it would not be considered spam. Still other users may actually desire to receive certain types of commercial email that would be considered spam by other users. Today's email filter and anti-spam systems often fail to provide for such conditions.

[06] Thus, it is desirable to improve detection of bulk and/or spam email.

SUMMARY OF THE INVENTION

[07] A preferred embodiment of the invention provides for detection of "bulk" email by comparing email messages to one another to determine if the email messages should be considered essentially the same message. If a category of "same" messages meets a predetermined criteria, such as the total number of messages in the category exceeding a predefined number, then the messages are considered bulk messages. Once a determination is made that email messages are bulk email messages then further analysis can be performed to determine if the bulk email is unwanted, or "spam." Depending on the spam determination further action can be taken, such as preventing the messages from delivery.

A determination of "sameness" for purposes of both bulk and spam classifications can use any number and type of evaluation modules. Each module can include one or more

rules, tests, processes, algorithms, or other functionality. For example, one type of module may be a word count of email message text. Another module can use a weighting factor based on groups of multiple words and their perceived meanings. In general, any type of module can be used. A preferred embodiment of the invention uses statistical analysis, such as Bayesian analysis, to measure the performance of different modules against a known standard, such as human manual matching. Modules that are performing worse than other modules can be valued less than modules having better performance. In this manner, a high degree of reliability can be achieved. To improve performance, if a message is determined to be the same as a previous message, the previous computations and results for that previous message can be re-used.

[08] In one embodiment, users are provided with options to customize or regulate bulk and spam classification and subsequent actions on how to handle the classified email messages. For example, a user can set parameters to select which modules are used, set the threshold number of "same" messages requirement for classifying a group of messages as bulk, set confidence limits for classifications, set filters based on number of matches of sameness engines, set "hold" times for incoming messages pending a determination of bulk classification, etc. Messages which are already determined to be bulk do not need to be held.

[09] One embodiment of the invention provides for a central server to distribute module performance and module related info data to different servers for use by the servers in sameness determinations. The central server handles the computation on module performance, create module related data, and can assist the servers in switching over from one module, or set of modules, to another. The switchover to a new module set can be performed over time at the direction of the central server to take into account changing characteristics of spam email, or to take into account changing tactics of spam senders.

[10] In one embodiment the invention provides an apparatus for classifying email messages, the apparatus comprising a processor for executing instructions included in a machine-readable medium, the machine-readable medium including one or more instructions for using a plurality of modules to determine a level of sameness of a particular email message with one or more prior email messages, wherein the level of sameness is derived for the particular email message from a weighting of the outputs of the modules; one or more instructions for determining a performance level for each of the modules; one or more instructions for comparing performance levels; one or more instructions for adjusting a weighting of at least one module in

response to comparing performance levels; and one or more instructions for using the level of sameness for the particular email message to classify the particular email message into a category.

[11] In another embodiment the invention provides a method for classifying
5 email messages, the method comprising using a plurality of modules to determine a level of sameness of a particular email message with one or more prior email messages, wherein the level of sameness is derived for the particular email message from a weighting of the outputs of the modules; determining a performance level for each of the modules; comparing performance
10 levels; adjusting a weighting of at least one module in response to comparing performance levels; and using the level of sameness for the particular email message to classify the particular email message into a category.

BRIEF DESCRIPTION OF THE DRAWINGS

15 Fig. 1 is a prior art diagram of an email network; and

Fig. 2 illustrates basic processing and information transfers according to a preferred embodiment of the present invention.

DETAILED DESCRIPTION OF THE INVENTION

20 [12] Fig. 1 illustrates general characteristics of an email system using a digital network such as the Internet. Although the invention is discussed primarily with respect to email transferred over the Internet, any suitable network, network topology, transmission protocols, sender-receiver devices and relationships, and other characteristics or properties of electronic devices, processes and transmission methods can be used. For example, features of the invention
25 can be employed on a smaller scale to local area networks (LANs), campus or corporate networks, home networks, etc.

[13] In Fig. 1, a sender at a client computer such as client computer 10 can create and send an email message to email server 20. Email server 20 transfers the email over Internet 30 to receiving email server 40. Receiving email server 40 transfers the email message
30 to the intended recipient computer 50. Note that Fig. 1 is intended to be a general depiction of

one possible configuration of hardware and resources to achieve email transfer. As is known in the art, many types of devices and arrangements for email transfer are possible.

[14] The Internet can include any number of servers, routers, switches and other devices through which the email information travels before reaching an intended recipient or other final destination (i.e., receiver). Although the invention is discussed in connection with a client/server architecture, such nomenclature is for convenience, only, as the roles and functions of any specific processor can be that of a client or server at different times. In some cases, an email server, or other type of server need not be used, or the server functions can be integrated with a client or other processing device.

[15] Any type of processing devices can be used to send and receive email. For example, portable computing devices such as a personal digital assistant (PDA), cell phone, laptop computer, or other devices can be employed. In general, the devices and manner of specific processing (including location and timing) are not critical to practicing important features of the present invention.

[16] A preferred embodiment of the invention uses an email server, or central server. The central server can be either of servers 20 or 40, or it can be another server through which an email message is transferred. For example, an Internet service provider (ISP) or other entity that provides email services may operate one or more email servers (not shown). In most embodiments it is desirable to have one or a few centralized points through which email traffic flows in order to be able to analyze and filter email to eliminate unwanted email, or spam.

[17] Fig. 2 illustrates basic processing and information transfers according to a preferred embodiment of the present invention.

[18] In Fig. 2, system 100 includes central server 120 for receiving incoming email traffic 105. Incoming email traffic is processed by identification (ID) modules 130. ID modules 130 can each be independent processes, devices or other types of functionality that are designed to evaluate the similarity between two or more email messages, or between an email message and other reference data. For example, a module can be a software process that determines a word count in the body of an email message. Another module can eliminate minor differences in grammar and language properties (e.g., eliminating the distinction between number, tense and person in grammar) and then subject the email text to a hash function. Another module type can use a database to assign values to words and compute an overall value

to an email message or other information in an email transmission. Another module type can attempt to deduce meanings or concepts conveyed in an email message. In general, any type of module that performs a comparison on an email message in an attempt to find a level or degree of “sameness” of the email message with other email messages or with a reference datum or data structure can be employed.

[19] Modules 130 are used to produce a value or values that indicate the degree of certainty that the module assigns to an email message to indicate whether the email message is matched to another email message or group of email messages (referred to as a category). In a preferred embodiment, a module can merely produce a “match” or “no match” with a message category. For example, there maybe tens or hundreds or more different categories to which a candidate email message may be compared. Each module can produce a conclusion as to whether the candidate email message is matched to one or more of the categories. Alternatively, the modules can produce a “degree of certainty” value (e.g., from 0 to 100) that indicates the degree to which the module ascertains that the message is the same as the messages in each category. Other types of module comparisons and types of module output are possible. Typically, the checking for sameness is performed against one or more messages, or other information representing message content, obtained from a source such as database 150.

[20] Filter 140 uses the results of the module comparisons to determine whether the email traffic is delivered to a user such as user 160. Messages that are delivered, or “passed,” are considered legitimate, or desired, messages. Messages that are undelivered are said to be blocked, or undesirable, messages. In a preferred embodiment, filter 140 performs a spam check on messages that are considered to be “bulk” messages. Bulk messages are those messages in a category that contains over X messages, where X is a predetermined number. Messages that are not bulk messages are automatically passed through to the intended recipient, or user. Bulk messages, on the other hand, are subjected to further scrutiny and checking to determine whether the messages are spam. Typically, this is a determination as to whether the message is of a commercial nature, i.e., is designed to benefit the sender monetarily or in some other fashion. The spam checking can be by any means as is known in the art.

[21] Ideally, a system would be able to tell with certainty whether a message is bulk or not. In such a case, all non-bulk messages can be considered as not spam. However, since the bulk detection is not perfect, a preferred embodiment uses the probability that email is

bulk combined with the probability that the email may be spam to derive the overall probability that email is bulk and spam.

[22] Performance evaluation 110 is a process that assesses the success or failure of specific modules to detect a match of email. For example, if a module has determined that an email message is the same as other email messages in category 1, and a benchmark evaluation reveals that the message is actually NOT the same as the email messages in category 1, the module is said to have failed at the matching process. Modules that fail more frequently than desired can be removed from use, or if they are still used less weight can be given to their conclusions. On the other hand, modules which are more successful than others can be weighted more heavily than other modules so that the successful modules have more influence in the overall decision of sameness than other modules. Completely useless or redundant modules can be removed. A module X is redundant if there exists another module Y which is right whenever module X is right. A partially useless module may still be deleted if the incremental value that it provides does not justify the cost of executing the module.

[23] A benchmark evaluation can be achieved manually, as where it is performed by a human editor. Such an approach can periodically have a human editor review categorized email messages and determine if the sorting is proper. Module decisions as to the category can then be compared against the human editor and modules can be rated accordingly. Inter-module checking can also be performed where if a module often decides contrary to other modules, especially when the other modules are weighted highly, then the contrary module can be de-rated, or given a lower weighting in the overall result. A preferred embodiment uses Bayesian analysis to achieve this "tuning" of the module selection and weighting.

[24] Module tuning and sameness checking can be applied to spam checking in addition to bulk detection. A preferred embodiment of the invention uses a process called the SpamRepute Engine that uses the "sameness" method. If one message is identified to be spam, then all messages which are identified as the "same" as that message are also classified as spam. If one message is identified to be legitimate, then all messages which are identified as the "same" as that message is also legitimate.

[25] Table I shows examples of different ID modules used to perform sameness checking.

Module Name	Description
Nilsimsa	Uses nilsimsa fingerprints
DCC	Uses DCC checksums of messages, utilizes dccproc to get checksums
WordsWB	Using the whole message body, this module creates a hash
WordsF2K	Using only the first 2K bytes of the message body, this module creates a hash
WordsHB	Using only the first half of the message body, this module creates a hash

TABLE I

[26] Table II shows performance measurements of the modules in Table I

(a) Module Name	(b) #says matches and is right	(c) #says matches and is wrong	(d) #says not matches and is right	(e) #says not matches and is wrong	(f) Bayesian Computed weight
Nilsimsa	4317	5679	11691	813	0.87
DCC	2439	0	17370	2691	0.88
WordsWB	4571	669	16701	559	0.96
WordsF2K	4560	750	16620	341	0.98
WordsHB	4825	836	16534	305	0.98

TABLE II

5

[27] In Table II, a determination as to whether the module is actually right or wrong is based on a human editor's decision as to whether the messages are the "same." One way to calculate a weight to assign to a module based on the results of Table II can be performed as follows (using the values denoted by the letters in parentheses in each column of Table II):

$$\text{ID Module sameness probability} = \frac{(b/(b+c))}{(b/(b+c) + e/(d+e))}$$

[28] The calculation above is computed for each module used in a group of modules (e.g., ID modules 130 of Fig. 2) to evaluate sameness. Then, the preferred embodiment uses the overall "sameness" probability derived by combining the results of all of the modules using a chi-squared probability scheme which is described, for example, at <http://mathworld.wolfram.com/Chi-SquaredDistribution.html>. There are many other alternative schemes to combine the module results and any suitable approach can be used.

[29] By assuming that messages with a computed 99% or greater combined sameness probability are the same and those with a 1% or less combined sameness probability are not the same, the message sample can be used to tune the weights of the individual ID modules. In order to start a comparison process to determine if an email message is spam, it may
 5 be necessary to make a determination about at least one instance of the email message. One way this can be done is by observation and flagging an email as spam by a human editor. Another method is to allow a user or group of users to “vote” on whether the email is spam. The votes can be cross referenced by the ID modules, weighted by a user’s reputation and tallied.

[30] Other ways to perform initial classification can include so-called “relay
 10 honeypots” and “proxy honeypots” to obtain information about an email message or email sender. Another approach can use sender authentication and permit messages from authenticated senders, or with approved digital signatures, to be passed as legitimate, non-spam, email.

[31] The Source Code Appendix includes examples of ID modules and a routine to manage the application of the modules to determine sameness among email messages.

[32] One feature of the present invention allows a central server to hold
 15 messages for a time before delivery. In this manner a group of messages can be compared for sameness before they are delivered to the intended recipients. This can be useful, for example, to make a bulk email determination since the threshold criterion for bulk email delivery might be tens of thousands of messages but an email spammer may only send out a few hundred messages
 20 at a time so that the bulk detection could not occur until after a reasonable time interval (e.g., 1 hour). If the non-bulk messages are not delayed, then all messages below the bulk threshold will always be treated as non-bulk. The time period for holding email messages pending delivery in order to detect bulk email can vary according to any of a number of factors. For example, a user may designate certain types of email (e.g., from known senders) to be delivered immediately.
 25 The number to be held could be proportional to the historical number of messages needed for actual bulk messages to exceed the bulk threshold. Email from senders who are not on the immediate delivery list can be held for a predetermined time period. The time period for holding can also be increased gradually as the number of matches for a category grows or based on the content of the message. Other ways to set or vary the hold time are possible.

[33] A contemplated embodiment of the invention envisions a master server
 30 performing the performance evaluation function shown in Fig. 2. The results of performance

evaluation can be sent to one or more other servers so that the other servers are not burdened with the potentially time and resource-consuming operation of evaluating the modules. The master server can make recommendations on module weighting and can phase in or out modules for use in bulk detection and/or spam detection. Module definitions and module related
5 databases are synchronized between the client and server.

[34] Another feature allows for users to set parameters that deal with tuning sameness engines for either bulk mail detection or spam detection. For example, a user may want to include or exclude certain modules, or to manually assign different weights to different modules. This can be useful, for example, where a user wants to receive certain types of email that might otherwise be blocked as spam. A user can be allowed to set any of a number of
10 parameters in the system, such as the threshold number of email messages in a category before the category is considered to be bulk email. Note that many different users can each have a different threshold number so that the same category of email messages can be considered bulk (or spam) or not on a user-by-user basis. Subsequent processing or detection of spam or other
15 characteristics of email can proceed, accordingly, on the same user-by-user basis. User selection or setting of parameters can be by typing in a numerical value, selection of a menu option, activation of a button or slider control, or by any other means as is known in the art. If a continuous function is converted into a discrete function, the confidence level determines the threshold by which the conversion is performed. For example, if the probability is 99 and the
20 confidence level is 98, then the result can be converted into a definite 100.

[35] Although the invention has been discussed with reference to specific embodiments thereof, these embodiments are illustrative, and not restrictive, of the invention. For example, although the invention is discussed primarily with respect to email messages, other
25 types of information exchange or transfer can similarly benefit. For example, so-called pop-up web pages or dialog boxes can be subjected to sameness testing and filtering. Items attached to email messages can be analyzed as can web pages, files, images, or other data items or structures.

[36] In the description herein, numerous specific details are provided, such as
30 examples of components and/or methods, to provide a thorough understanding of embodiments of the present invention. One skilled in the relevant art will recognize, however, that an

embodiment of the invention can be practiced without one or more of the specific details, or with other apparatus, systems, assemblies, methods, components, materials, parts, and/or the like. In other instances, well-known structures, materials, or operations are not specifically shown or described in detail to avoid obscuring aspects of embodiments of the present invention.

5 [37] A “machine-readable medium” or “computer-readable medium” for purposes of embodiments of the present invention may be any medium or transmission that can contain, store, communicate, propagate, or transport the program for use by or in connection with the instruction execution system, apparatus, system or device. The computer readable carrier can be, by way of example only but not by limitation, an electronic, magnetic, optical,
10 electromagnetic, infrared, or semiconductor system, apparatus, system, device, propagation medium, or computer memory.

 [38] A “processor” or “process” includes any human, hardware and/or software system, mechanism or component that processes data, signals or other information. A processor can include a system with a general-purpose central processing unit, multiple processing units,
15 dedicated circuitry for achieving functionality, or other systems. Processing need not be limited to a geographic location, or have temporal limitations. For example, a processor can perform its functions in “real time,” “offline,” in a “batch mode,” etc. Portions of processing can be performed at different times and at different locations, by different (or the same) processing systems.

20 [39] Reference throughout this specification to “one embodiment”, “an embodiment”, or “a specific embodiment” means that a particular feature, structure, or characteristic described in connection with the embodiment is included in at least one embodiment of the present invention and not necessarily in all embodiments. Thus, respective appearances of the phrases “in one embodiment”, “in an embodiment”, or “in a specific
25 embodiment” in various places throughout this specification are not necessarily referring to the same embodiment. Furthermore, the particular features, structures, or characteristics of any specific embodiment of the present invention may be combined in any suitable manner with one or more other embodiments. It is to be understood that other variations and modifications of the embodiments of the present invention described and illustrated herein are possible in light of the
30 teachings herein and are to be considered as part of the spirit and scope of the present invention.

[40] Embodiments of the invention may be implemented by using a programmed general purpose digital computer, by using application specific integrated circuits, programmable logic devices, field programmable gate arrays, optical, chemical, biological, quantum or nanoengineered systems, components and mechanisms may be used. In general, the functions of the present invention can be achieved by any means as is known in the art. Distributed, or networked systems, components and circuits can be used. Communication, or transfer, of data may be wired, wireless, or by any other means.

[41] It will also be appreciated that one or more of the elements depicted in the drawings/figures can also be implemented in a more separated or integrated manner, or even removed or rendered as inoperable in certain cases, as is useful in accordance with a particular application. It is also within the spirit and scope of the present invention to implement a program or code that can be stored in a machine-readable medium to permit a computer to perform any of the methods described above.

[42] Additionally, any signal arrows in the drawings/Figures should be considered only as exemplary, and not limiting, unless otherwise specifically noted. Furthermore, the term “or” as used herein is generally intended to mean “and/or” unless otherwise indicated. Combinations of components or steps will also be considered as being noted, where terminology is foreseen as rendering the ability to separate or combine is unclear.

[43] As used in the description herein and throughout the claims that follow, “a”, “an”, and “the” includes plural references unless the context clearly dictates otherwise. Also, as used in the description herein and throughout the claims that follow, the meaning of “in” includes “in” and “on” unless the context clearly dictates otherwise.

[44] The foregoing description of illustrated embodiments of the present invention, including what is described in the Abstract, is not intended to be exhaustive or to limit the invention to the precise forms disclosed herein. While specific embodiments of, and examples for, the invention are described herein for illustrative purposes only, various equivalent modifications are possible within the spirit and scope of the present invention, as those skilled in the relevant art will recognize and appreciate. As indicated, these modifications may be made to the present invention in light of the foregoing description of illustrated embodiments of the present invention and are to be included within the spirit and scope of the present invention.

[45] Thus, while the present invention has been described herein with reference to particular embodiments thereof, a latitude of modification, various changes and substitutions are intended in the foregoing disclosures, and it will be appreciated that in some instances some features of embodiments of the invention will be employed without a corresponding use of other features without departing from the scope and spirit of the invention as set forth. Therefore, many modifications may be made to adapt a particular situation or material to the essential scope and spirit of the present invention. It is intended that the invention not be limited to the particular terms used in the following claims and/or to the particular embodiment disclosed as the best mode contemplated for carrying out this invention, but that the invention will include any and all embodiments and equivalents falling within the scope of the appended claims.